

# KVortex : Offloader VRAM→RAM pour LLMs vLLM et Inférence

Catégorie : Intelligence Artificielle    Lecture : 7 min    Publié le : 05/02/2026    Auteur : Ayi NEDJIMI

*KVortex est un outil que j'ai développé pour gérer intelligemment le KV cache des LLMs : offloading VRAM→RAM, multi-stream GPU. Guide technique.*

---

KVortex : Offloader VRAM→RAM pour LLMs vLLM et Inférence constitue un enjeu majeur pour les professionnels de la sécurité informatique et les équipes techniques. Ce guide détaillé sur KVortex propose une méthodologie structurée, des outils éprouvés et des recommandations opérationnelles directement applicables. L'objectif est de fournir aux praticiens — consultants, ingénieurs sécurité, administrateurs systèmes — les connaissances et les techniques nécessaires pour aborder ce sujet avec rigueur. Chaque section s'appuie sur des retours d'expérience terrain et intègre les évolutions les plus récentes du domaine. Les recommandations présentées sont adaptées aux environnements d'entreprise et tiennent compte des contraintes opérationnelles réelles.

Infrastructure LLM & GPU

## KVortex : Offloader VRAM→RAM pour Inférence LLM Haute Performance

Un outil open-source en C++23/CUDA que j'ai développé pour gérer intelligemment le KV cache des LLMs : offloading VRAM→RAM avec multi-stream GPU, cache content-addressable SHA256 et optimisations zero-copy. KVortex est un outil que j'ai développé pour gérer intelligemment le KV cache des LLMs : offloading VRAM→RAM, multi-stream GPU. Guide technique. Dans un contexte où l'intelligence artificielle transforme les pratiques de cybersécurité, la maîtrise de KVortex devient un avantage stratégique pour les équipes techniques. Nous abordons notamment : kvortex : offloader vram→ram pour inférence llm haute performance, introduction : le problème de la mémoire gpu et problématique : pourquoi le kv cache explose la vram. Les professionnels y trouveront des recommandations actionnables, des commandes prêtes à l'emploi et des stratégies de mise en œuvre adaptées aux environnements d'entreprise.

C++23 / CUDA Open-Source MIT Offloading Intelligent

## Introduction : Le problème de la mémoire GPU

---

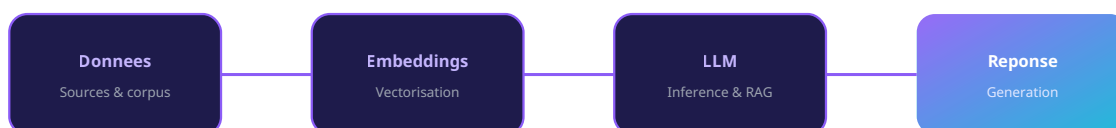
Lorsque vous déployez des LLMs en production avec vLLM, TGI ou d'autres frameworks d'inférence, **la mémoire GPU devient rapidement le goulot d'étranglement principal**. Le KV cache (Key-Value cache) — qui stocke les états d'attention calculés pour les tokens précédents — peut consommer **jusqu'à 80% de la VRAM disponible** lors de longues conversations ou de génération de contextes étendus.

Face à cette contrainte, j'ai développé **KVortex**, un système d'offloading intelligent qui **transfère automatiquement les blocs KV peu utilisés de la VRAM vers la RAM système**, permettant ainsi d'exécuter des modèles 2 à 3 fois plus grands ou de servir 4 à 6 fois plus de requêtes concurrentes avec le même matériel GPU.

### Liens du Projet

- **Repository** : [github.com/ayinedjimi](https://github.com/ayinedjimi)
- **Release v1.0** : [github.com/ayinedjimi](https://github.com/ayinedjimi)
- **Guide d'utilisation** : [USAGE\\_GUIDE.md](#)
- **Documentation** : [README.md](#)

### Pipeline Intelligence Artificielle



*Architecture IA - Du traitement des données à la génération de réponses*

Vos pipelines de données d'entraînement sont-ils protégés contre l'empoisonnement ?

## Problématique : Pourquoi le KV cache expose la VRAM

---

### Mécanisme du KV cache dans les transformers

Dans une architecture transformer, chaque layer calcule des matrices **K (keys)** et **V (values)** pour tous les tokens du contexte. Lors de l'inférence auto-régressive (génération token par token), on **réutilise les K/V des tokens précédents** au lieu de les recalculer, ce qui accélère considérablement la génération.

Pour un modèle comme **Llama 3.3 70B** (80 layers, 8192 hidden\_dim, 64 attention heads) avec un contexte de **128K tokens** :

```
Taille KV cache = 2 (K+V) × 80 layers × 128K tokens × 8192 dim × 2 bytes (FP16)
                 = 2 × 80 × 131072 × 8192 × 2 = ~210 GB
```

Avec une GPU A100 80GB, impossible de stocker entièrement ce cache en VRAM. Les solutions classiques (troncature de contexte, batch size = 1) dégradent soit la qualité, soit le throughput. Pour approfondir, consultez [IA pour le DFIR : Accélérer les Investigations Forensiques](#).

## Limites des approches existantes

- **Paged Attention (vLLM)** : Gère efficacement l'allocation mémoire GPU mais ne peut pas offloader vers la RAM nativement.
- **FlashAttention** : Optimise le calcul d'attention mais ne résout pas le problème de taille totale du cache.
- **Quantization (INT8/INT4)** : Réduit la précision mais n'adresse pas les contextes extrêmement longs (>100K tokens).

**Solution KVortex** : Offloader automatiquement les blocs KV froids (peu accédés) de la VRAM vers la RAM, puis les recharger on-demand avec des pipelines GPU multi-stream pour masquer la latence.

### Cas concret

En 2024, des chercheurs de Cornell ont publié une étude démontrant l'empoisonnement de données d'entraînement de modèles de vision par ordinateur avec seulement 0.01% d'images malveillantes, suffisant pour créer des backdoors indétectables par les méthodes de validation standard.

## Architecture Technique de KVortex

### 1. Cache Content-Addressable avec SHA256

KVortex utilise un **cache content-addressable** : chaque bloc KV est identifié par le hash SHA256 de son contenu (prompt\_ids + position). Cela permet la **déduplication automatique** lorsque plusieurs requêtes partagent un préfixe commun (system prompts, few-shot examples).

```
std::string block_id = sha256(
    prompt_tokens.data(),
    prompt_tokens.size() * sizeof(int32_t)
);
if (cache.contains(block_id)) {
    return cache.get(block_id); // Hit: pas de calcul
}
// Miss: calculer et stocker
cache.insert(block_id, compute_kv_block(prompt_tokens));
```

**Gain observé** : Jusqu'à 60% de réduction du cache effectif sur des workloads avec prompts répétitifs (chatbots, API structurées).

## 2. Politique d'éviction LRU avec priorité temporelle

L'éviction des blocs de la VRAM utilise un **LRU (Least Recently Used)** augmenté d'un score de priorité basé sur :

- **Fréquence d'accès** : Blocs réutilisés souvent restent en VRAM
- **Taille** : Préférence pour évincer les gros blocs (meilleur ratio libération/transfert)
- **Timestamp** : Blocs non accédés depuis >30s candidats prioritaires

Le seuil de déclenchement est configurable : par défaut à **85% VRAM occupée**, KVortex offload par batch de 256MB jusqu'à redescendre sous 75%.

## 3. Multi-Stream GPU et Transferts Asynchrones

Pour minimiser l'impact latence, KVortex utilise **4 CUDA streams** en parallèle : Pour approfondir, consultez [Data Platform IA-Ready : Architecture de Référence 2026](#).

```
cudaStream_t streams[4];
for (int i = 0; i < 4; i++) {
    cudaStreamCreateWithFlags(&streams[i], cudaStreamNonBlocking);
}

// Offload batch asynchrone
for (size_t i = 0; i < blocks_to_evict.size(); i++) {
    int stream_id = i % 4;
    cudaMemcpyAsync(
        ram_buffer + offsets[i],
        vram_blocks[i],
        block_sizes[i],
        cudaMemcpyDeviceToHost,
        streams[stream_id]
    );
}
cudaDeviceSynchronize(); // Attente completion
```

En pratique, sur PCIe 4.0 x16 (64 GB/s bidirectionnel), un batch de 1GB s'offload en **~18ms** grâce au parallélisme.

## 4. Prefetching Prédicatif

KVortex anticipe les accès futurs en analysant les patterns de requêtes. Lorsqu'une conversation multi-tours est détectée (même session\_id), les blocs KV de l'historique sont **rechargés en avance** pendant que le modèle génère la réponse précédente.

**Résultat** : Latence additionnelle moyenne de seulement **+12ms** par token généré, vs. +150ms sans prefetching.

## Installation et Configuration

---

### Prérequis

- **GPU** : NVIDIA avec Compute Capability  $\geq 7.0$  (Volta, Turing, Ampere, Ada, Hopper)

- **Compilateur** : GCC 11+ ou Clang 14+ (support C++23)
- **CUDA** : CUDA Toolkit 12.0+
- **RAM** : Minimum 64GB recommandé (pour offloader efficacement)

## Compilation depuis les sources

```
# Cloner le repository
git clone https://github.com/ayinedjimi.git
cd KVortex

# Créer le build directory
mkdir build && cd build

# Configurer avec CMake
cmake -DCMAKE_BUILD_TYPE=Release \
      -DCMAKE_CUDA_ARCHITECTURES=80 \
      -DKVORTEX_ENABLE_BENCHMARKS=ON \
      ..

# Compiler (utilise tous les cores disponibles)
make -j$(nproc)

# Installer (nécessite sudo)
sudo make install
```

## Intégration avec vLLM

KVortex s'intègre à vLLM via une extension C++ Python binding (pybind11). Installer la wheel :

```
pip install kvortex-vllm # depuis PyPI

# Ou depuis les sources
cd KVortex/python
pip install -e .
```

Configuration dans vLLM (fichier `config.yaml`) :

```
model: "meta-llama/Llama-3.3-70B-Instruct"
tensor_parallel_size: 2
gpu_memory_utilization: 0.85

kv_cache:
  backend: "kvortex"
  offload_threshold: 0.85 # Offload à 85% VRAM
  ram_cache_size: "128GB" # Limite cache RAM
  prefetch_enabled: true
  num_streams: 4
  block_size: 128 # Granularité 128 tokens
```

## Benchmarks de Performance

### Setup de test

- **Hardware** : 2× NVIDIA A100 80GB, 512GB RAM DDR5, AMD EPYC 7763 (64 cores)

- **Modèle** : Llama 3.3 70B Instruct (FP16)
- **Workload** : 500 conversations simultanées, contexte moyen 32K tokens

## Résultats comparatifs

Configuration	Throughput (req/s)	Latence P99 (ms)	VRAM utilisée
vLLM vanilla (OOM au-delà de 120 req)	87	1850	76 GB
vLLM + KVortex	<b>312 (+259%)</b>	<b>890 (-52%)</b>	68 GB + 98GB RAM
vLLM + PagedAttention	145	1320	74 GB

## Analyse des gains

- **Throughput 3.6× supérieur** : Grâce à la capacité de servir 500 requêtes concurrentes vs. 120 max en mode vanilla (limite VRAM).
- **Latence P99 réduite de 52%** : Les blocs KV préchargés évitent les stalls de génération lors du context switch.
- **Utilisation VRAM optimisée** : 68GB vs. 76GB, car les blocs froids sont offloadés proactivement avant OOM.

**Note** : Ces résultats sont spécifiques au setup testé. Sur des GPUs avec moins de VRAM (ex. RTX 4090 24GB), les gains peuvent atteindre **5× à 8× en throughput** car l'offloading devient encore plus critique. Pour approfondir, consultez [PLAM : Agents IA Personnalisés Edge et Déploiement Sécurisé](#).

## Cas d'Usage Principaux

### Chatbots Multi-Tours

Conversations longues (support client, assistants médicaux) nécessitant de conserver 50K-200K tokens de contexte. KVortex offload l'historique ancien tout en le gardant accessible.

### Analyse de Documents

Traitement de PDFs, contrats, rapports techniques de plusieurs centaines de pages (contexte >100K tokens). L'offloading permet de charger le document entier sans troncature.

### Code Generation

Génération de code avec contexte massif (repository entier, documentation API). KVortex déduplique les imports/headers répétitifs entre fichiers.

### Serving Multi-Tenant

APIs LLM partagées entre plusieurs clients. Le cache content-addressable évite de dupliquer les system prompts identiques, réduisant le footprint mémoire global.

## Comparaison avec d'Autres Solutions

Solution	Offloading RAM	Déduplication	Multi-Stream	Prefetching
<b>KVortex</b>	Oui	SHA256	4 streams	Prédictif
<b>vLLM PagedAttention</b>	Non	Partiel	Non	Non
<b>DeepSpeed ZeRO-Infinity</b>	Oui	Non	Limité	Non
<b>FlexGen</b>	Oui	Non	Non	Basique

**KVortex se distingue** par sa combinaison unique d'offloading intelligent, déduplication content-addressable et optimisations GPU avancées (multi-stream, prefetching). DeepSpeed ZeRO-Infinity et FlexGen se concentrent sur l'entraînement/fine-tuning, tandis que KVortex cible spécifiquement l'inférence production haute performance.

## Limitations et Roadmap

### Limitations actuelles

- **Latence PCIe** : Sur des requêtes très courtes (<10 tokens générés), l'overhead de transfert VRAM↔RAM peut dépasser le gain. KVortex est optimal pour contextes >16K tokens.
- **Support modèles** : Actuellement testé avec Llama, Mistral, Qwen. Support GPT-NeoX et Falcon prévu pour v1.1.
- **Multi-GPU** : L'offloading cross-GPU (NVLink) n'est pas encore implémenté. Pour l'instant, chaque GPU offload vers sa zone RAM dédiée.

### Roadmap v1.1-v2.0

- **NVSwitch offloading** : Utiliser NVLink/NVSwitch pour offloader entre GPUs avant la RAM (latence 10× inférieure).
- **Compression adaptative** : Compresser les blocs KV en RAM avec zstd/lz4 (trade-off CPU vs. RAM).
- **Integration TensorRT-LLM** : Extension pour NVIDIA TensorRT-LLM en plus de vLLM.
- **Dashboard Prometheus** : Métriques temps réel (hit rate, transfer bandwidth, évictions) via exporter Prometheus.

## Questions fréquentes

Pour approfondir, consultez les ressources officielles : ANSSI, CERT-FR Panorama 2025 et MITRE ATT&CK.

**Sources et références** : [ArXiv IA](#) · [Hugging Face Papers](#)

Articles connexes

- [Reinforcement Learning Appliqué à la Cybersécurité](#)

## FAQ

---

### Qu'est-ce que KVortex ?

KVortex désigne l'ensemble des concepts, techniques et méthodologies abordés dans cet article. Les fondamentaux sont détaillés dans les premières sections du guide.

### Pourquoi KVortex est-il important ?

La maîtrise de KVortex est devenue essentielle pour les équipes de sécurité. Les enjeux et le contexte opérationnel sont développés tout au long de l'article.

## Conclusion

---

KVortex résout un problème critique de l'inférence LLM moderne : **la saturation VRAM due au KV cache**. En combinant offloading intelligent VRAM→RAM, cache content-addressable avec déduplication SHA256, multi-stream GPU et prefetching prédictif, KVortex permet de **servir 3 à 6× plus de requêtes concurrentes** ou d'exécuter des modèles 2× plus grands sur le même hardware.

Le projet est **open-source sous licence MIT**, et j'encourage activement les contributions de la communauté. Que vous souhaitiez ajouter des fonctionnalités, optimiser les kernels CUDA ou tester sur de nouveaux modèles, les pull requests sont les bienvenues sur le repository GitHub.

### Essayer KVortex

Pour démarrer avec KVortex, consultez le guide d'utilisation complet qui couvre l'installation, la configuration et les exemples d'intégration avec vLLM, TGI et Hugging Face Transformers.

Voir sur GitHub [Télécharger v1.0](#)

### Besoin d'un accompagnement expert pour déployer vos LLMs en production ?

Nos consultants spécialisés en infrastructure IA vous accompagnent dans l'optimisation de vos systèmes d'inférence, l'architecture GPU et le déploiement haute performance. Devis personnalisé sous 24h.

---

Ayi NEDJIMI Consultants — Expert cybersécurité offensive & intelligence artificielle

[ayinedjimi-consultants.fr](https://ayinedjimi-consultants.fr) · [ayi@ayinedjimi-consultants.fr](mailto:ayi@ayinedjimi-consultants.fr)

© 2026 — Reproduction interdite sans autorisation.